

# CUNY SPS DATA 621 - CTG5 - HW4

*Betsy Rosalen, Gabrielle Bartomeo, Jeremy O'Brien, Lidiia Tronina, Rose Koh*

*April 24th, 2019*

## Contents

<b>1</b>	<b>DATA EXPLORATION</b>	<b>2</b>
1.1	Summary Statistics . . . . .	3
1.2	Linearity . . . . .	9
1.3	Missing Data . . . . .	12
<b>2</b>	<b>DATA PREPARATION</b>	<b>13</b>
2.1	Missing Values . . . . .	13
<b>3</b>	<b>BUILD MODELS</b>	<b>16</b>
3.1	Classification Models: Models 1, 2, 3, 4 . . . . .	16
3.2	Regression Model: Models 5, 6 . . . . .	17
<b>4</b>	<b>SELECT MODELS</b>	<b>18</b>
<b>5</b>	<b>Appendix</b>	<b>18</b>

Table 1: Data Dictionary

VARIABLE	DEFINITION	TYPE
TARGET_FLAG	car crash = 1, no car crash = 0	binary categorical response
TARGET_AMT	car crash cost = >0, no car crash = 0	continuous numerical response
AGE	driver's age - very young/old tend to be risky	continuous numerical predictor
BLUEBOOK	\$ value of vehicle	continuous numerical predictor
CAR_AGE	age of vehicle	continuous numerical predictor
CAR_TYPE	type of car (6types)	categorical predictor
CAR_USE	usage of car (commercial/private)	binary categorical predictor
CLM_FREQ	number of claims past 5 years	discrete numerical predictor
EDUCATION	max education level (5types)	categorical predictor
HOMEKIDS	number of children at home	discrete numerical predictor
HOME_VAL	\$ home value - home owners tend to drive more responsibly	continuous numerical predictor
INCOME	\$ income - rich people tend to get into fewer crashes	continuous numerical predictor
JOB	job category (8types, 1missing) - white collar tend to be safer	categorical predictor
KIDSDRV	number of driving children - teenagers more likely to crash	discrete numerical predictor
MSTATUS	marital status - married people drive more safely	categorical predictor
MVR PTS	number of traffic tickets	continuous numerical predictor
OLDCLAIM	\$ total claims in the past 5 years	continuous numerical predictor
PARENT1	single parent	binary categorical predictor
RED_CAR	a red car	binary categorical predictor
REVOKE	license revoked (past 7 years) - more risky driver	binary categorical predictor
SEX	gender - woman may have less crashes than man	binary categorical predictor
TIF	time in force - number of years being customer	continuous numerical predictor
TRAVTIME	distance to work	continuous numerical predictor
URBANCITY	urban/rural	binary categorical predictor
YOJ	years on job - the longer they stay more safe	continuous numerical predictor

## 1 DATA EXPLORATION

In the pursuit of determining relationships between car crashes, their costs, and factors that may play a role into each, a dataset containing 8,161 observations with 25 variables was explored, analyzed, and modeled. This data came from an auto insurance company with each observation representing one of their customers. Of the 25 variables, two were target variables (car crashes and car costs), and the other 23 were predictors. TARGET\_FLAG is a binary variable where a value of 1 indicates that the customer has made a claim related to a car crash and a value of 0 indicates they have not. The other target variable, TARGET\_AMT, is a continuous numerical variable whose value is the payout amount of a claim, if any. The remaining variables are split in their categorization; 13 are categorical and 10 are numerical.

This data was utilized to compose and evaluate several types of models with the following features:

- Logistic classification models that aim to predict the probability that a person crashes their car; and,
- Multiple linear regression models that aim to predict the amount of money it will cost if the person does crash their car.

The intended use case for these models is actuarial in nature: specifically, to calculate insurance rates commensurate with policyholders' (or policy applicants') potential risk levels based on attributes such as income, age, distance to work, tenure as customers, so on and so forth.

Inspection of the target variables reveals that where TARGET\_FLAG has values of 0 (i.e., no claim), TARGET\_AMT also has values of 0 (i.e., no payout), which is logically consistent.

Table 2: (#tab:t2.1)Summary statistics

	n	min	mean	median	max	sd
TARGET_AMT	8161	0	1504.3	0	107586	4704.0
AGE	8155	16	44.8	45	81	8.6
YOJ	7707	0	10.5	11	23	4.1
INCOME	7716	0	61898.1	54028	367030	47572.7
HOME_VAL	7697	0	154867.3	161160	885282	129123.8
TRAVTIME	8161	5	33.5	33	142	15.9
BLUEBOOK	8161	1500	15709.9	14440	69740	8419.7
TIF	8161	1	5.3	4	25	4.2
OLDCLAIM	8161	0	4037.1	0	57037	8777.1
MVR_PTS	8161	0	1.7	1	13	2.1
CAR_AGE	7651	0	8.3	8	28	5.7

Table 3: (#tab:t2.2)Summary statistics for Categorical Variables

EDUCATION	JOB	CAR_TYPE	KIDSDRIV	HOMEKIDS	CLM_FREQ
<High School:1203	Blue Collar :1825	Minivan :2145	0:7180	0:5289	0:5009
Bachelors :2242	Clerical :1271	Panel Truck: 676	1: 636	1: 902	1: 997
Masters :1658	Professional:1117	Pickup :1389	2: 279	2:1118	2:1171
PhD : 728	Manager : 988	Sports Car : 907	3: 62	3: 674	3: 776
High School :2330	Lawyer : 835	Van : 750	4: 4	4: 164	4: 190
NA	Student : 712	SUV :2294	NA	5: 14	5: 18
NA	(Other) :1413	NA	NA	NA	NA

## 1.1 Summary Statistics

Continuous and categorical variables were summarized separately for the sake of clarity.

EDUCATION, JOB, CAR\_TYPE, KIDSDRIV, HOMEKIDS, and CLM\_FREQ each comprise multiple categories. On the other hand, PARENT1, SEX, MSTATUS, CAR\_USE, RED\_CAR, REVOKED, URBANICITY are all binaries.

### 1.1.1 Variable Descriptions

#### 1.1.1.1 KIDSDRIV

KIDSDRIV is a categorical predictor with values ranging from 0 to 4. It shows heavy skew, with most cars having no kid drivers (value of 0). Judging from the distribution, it appears that having a kid driver results in higher probability of making a claim.

#### 1.1.1.2 AGE

AGE presents driver's age and shows a normal distribution centered around 45 years. Looking at the boxplot of age below, there does not appear to be a difference in the distribution between whether a claim is made or not. Accordingly, this AGE may not be helpful in determining the probability of making a claim.

#### 1.1.1.3 HOMEKIDS

HOMEKIDS is a predictor describing number of children at home ranging from 0 to 5.

Table 4: (#tab:t2.2)Summary statistics for Binary Categorical Variables

PARENT1	SEX	MSTATUS	CAR_USE	RED_CAR	REVOKED	URBANICITY
No :7084	M:3786	Yes:4894	Commercial:3029	no :5783	No :7161	Urban:6492
Yes:1077	F:4375	No :3267	Private :5132	yes:2378	Yes:1000	Rural:1669

#### 1.1.1.4 YOJ

YOJ is a predictor describing years on job. People who stay at a job for a longer time are believed to be safer drivers. Apart from those who are unemployed (values of 0), YOJ seems to show a normal distribution.

#### 1.1.1.5 INCOME

INCOME is a heavily skewed predictor variable, suggesting that outliers should be treated for modelling.

#### 1.1.1.6 HOME\_VAL

HOME\_VAL is a home value predictor variable. In theory, home owners tend to drive more responsibly. The difference between owners and renters (values of 0) is visible in the summary statistics graph.

#### 1.1.1.7 TRAVTIME

TRAVTIME is a predictor variable describing the distance to work. Long drives to work would suggest greater risk of an accident and claim. However, its graph shows a fairly normal distribution, such that this variable may not be helpful in determining the probability of making a claim.

#### 1.1.1.8 BLUEBOOK

BBLUEBOOK is a predictor variable describing the value of the car. The boxplot demonstrates that the lower value of the car, the higher chances of making a claim. It is conceivable that higher-priced cars are driven more carefully.

#### 1.1.1.9 TIF

TIF describes how long the customer has been with the insurance company. Plots reveal that the longer the tenure of a policyholder, the lower the likelihood of a claim - i.e. safe drivers tend to remain so.

#### 1.1.1.10 OLDCLAIM

OLDCLAIM is a predictor describing the value of claims made in the past 5 years. It is very heavily skewed as most policyholders do not make claims.

#### 1.1.1.11 CLM\_FREQ

CLM\_FREQ is a predictor that describes the frequency of claims in the past 5 years. It suggests that those who have made a claim in the past 5 years are more likely to make another claim.

#### 1.1.1.12 MVR\_PTS

MVR\_PTS is a predictor that describes motor vehicle record points. The rationale is that more traffic tickets suggests less safe driving and a higher likelihood of claims. It appears to be a highly significant variable as seen in boxplots.

### **1.1.1.13 CAR\_AGE**

CAR\_AGE describes the age of the policyholder's vehicle. One value is -3, which must be an error - this is corrected to 0.

### **1.1.1.14 PARENT1**

PARENT1 indicates whether a policyholder is a single parent. This variable has been factorized and relabeled as NumParents to describe the number of parents.

### **1.1.1.15 SEX**

SEX describes the gender of the driver. This variable has been factorized and relabeled as MALE, for which males receive a value of 1 and females a value of 0. It does not appear to be significant variable in the boxplots below.

### **1.1.1.16 MSTATUS**

MSTATUS describes the marital status of the policyholder. The rationale is that married people drive more safely. This variable has been factorized and relabeled as Single, for which married policyholders receive a value of 0 and unmarried a value of 1.

### **1.1.1.17 EDUCATION**

EDUCATION describes the education level of the driver. This variable is factorized. It may be correlated with INCOME.

### **1.1.1.18 JOB**

JOB describes the type of job the driver has. This variable is factorized. It may be correlated with INCOME. In theory policyholders with white collar jobs tend to drive more safely.

### **1.1.1.19 CAR\_TYPE**

CAR\_TYPE describes type of car. This variable is factorized.

### **1.1.1.20 CAR\_USE**

CAR\_USE describes how the vehicle is used. Commercial vehicles are driven more and may have an elevated probability of accidents and claims. This variable is factorized and relabeled as Commercial, for which a value of 0 means private use and a value of 1 means commercial use.

### **1.1.1.21 RED\_CAR**

RED\_CAR indicates whether the color of the vehicle is red. Red vehicles, especially sports cars, are associated with riskier driving and likelihood of claims. This variable is factorized.

### **1.1.1.22 REVOKED**

REVOKED describes whether a policyholders license has been revoked in the past 7 years. License revocation is associated with riskier driving. This variable is factorized. The boxplot reveals that policyholders who previously lost their license are more likely to file claims.

### 1.1.1.23 URBANICITY

URBANICITY describes whether driver lives in an urban area or a rural area. This variable has been factorized and relabeled as URBAN, for which a value of 0 means rural and a value of 1 means urban.

## 1.1.2 Summary Statistics Graphs

[GAB: This sentence needs rephrasing and a supporting chart... if the supporting chart is the next chart, then this needs to be/should be moved] Examining the dispersion of claims between variables, it looks like likelihoods are higher for drivers who are male, urban, blue collar, unmarried, or parents; as well as for those with commercial vehicles or a revoked license.

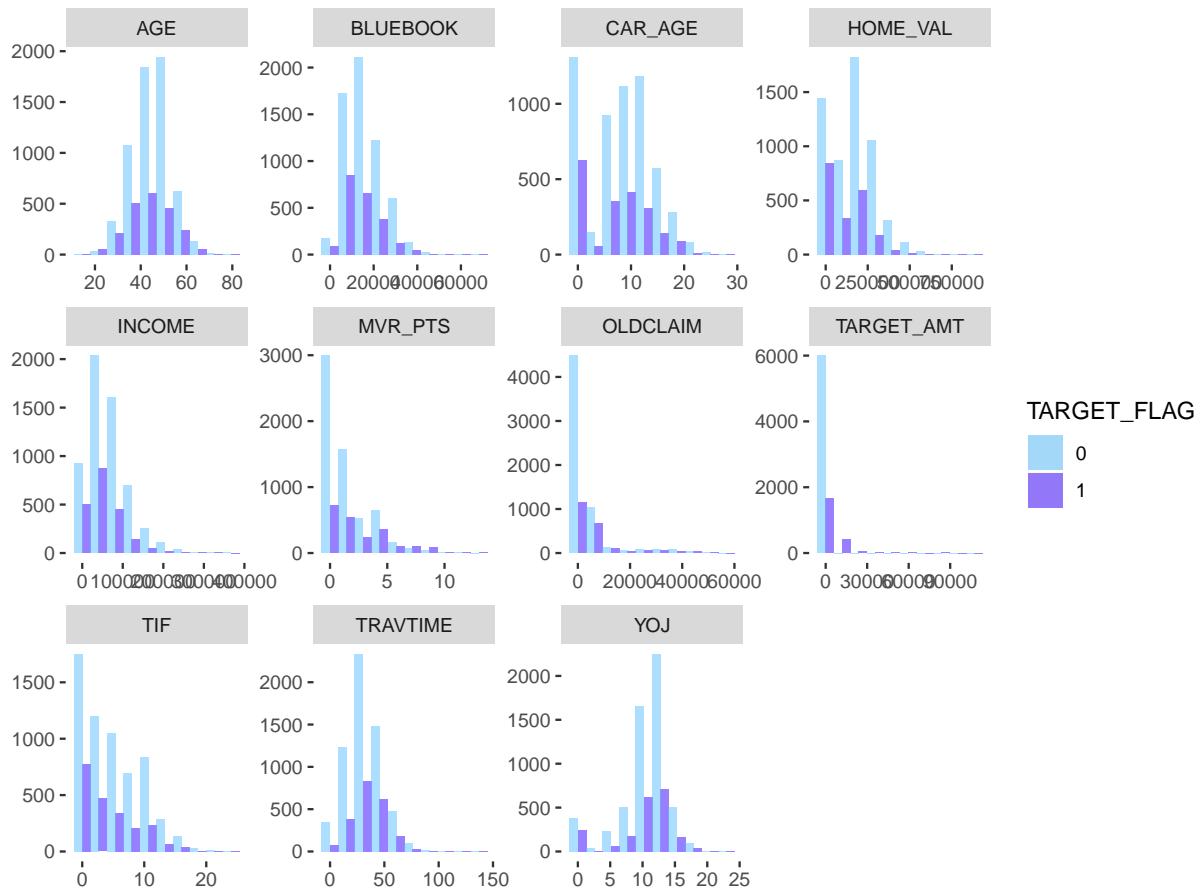


Figure 1: Numeric Data Distributions as a Function of TARGET\_FLAG

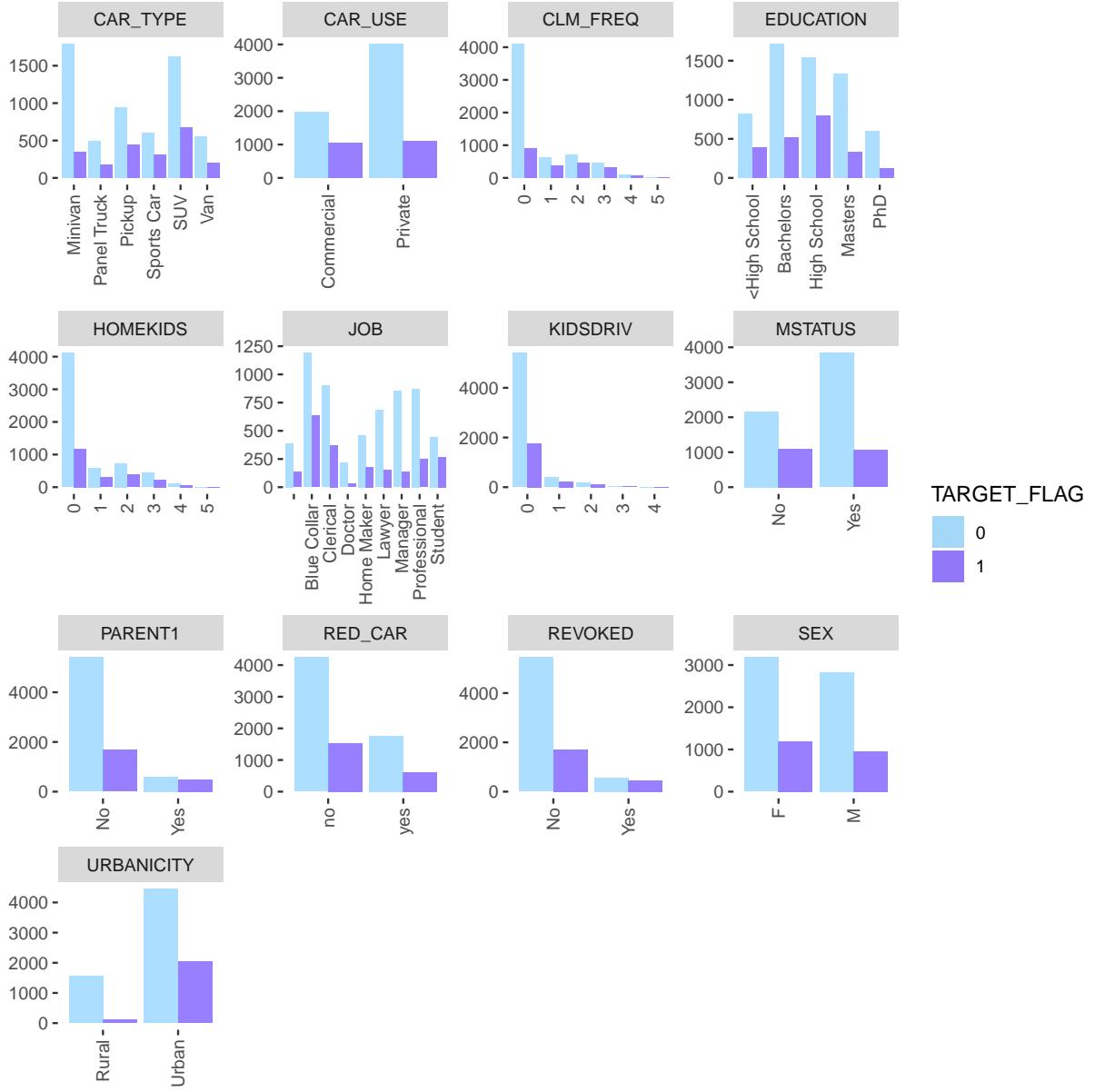


Figure 2: Categorical Data Distributions as a Function of `TARGET_FLAG`

The scale of the continuous variables' distributions are considerably different and difficult to visualize together.

Scaling the distribution based on the standard deviation reveals that outliers are very abundant for the continuous variables `OLDCLAIM`, `INCOME`, `TRAV_TIME`, `BLUEBOOK`, and to a lesser extent `HOME_VAL` and `TIF`. The variables that appear to have the most outliers are `OLDCLAIM`, `BLUEBOOK`, `TRAVTIME`, and `INCOME`. All of the variables show varying levels of skew save `YOJ` and `AGE` which appear the most normally distributed.

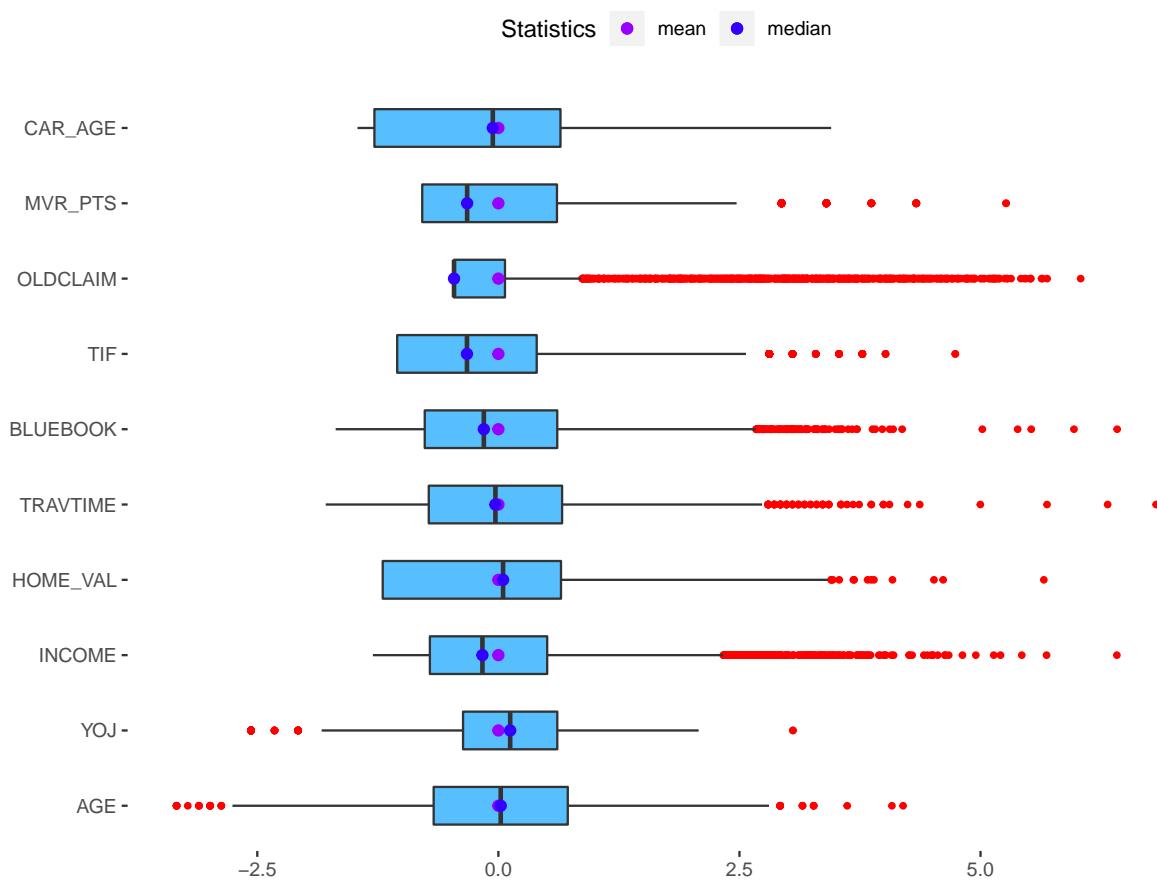


Figure 3: Scaled Boxplots

[JO: WEREN'T WE GO TO TOSS THE NEXT CHART?]

[GAB: I DON'T KNOW, WERE WE? NO TEXT UNTIL I KNOW!]

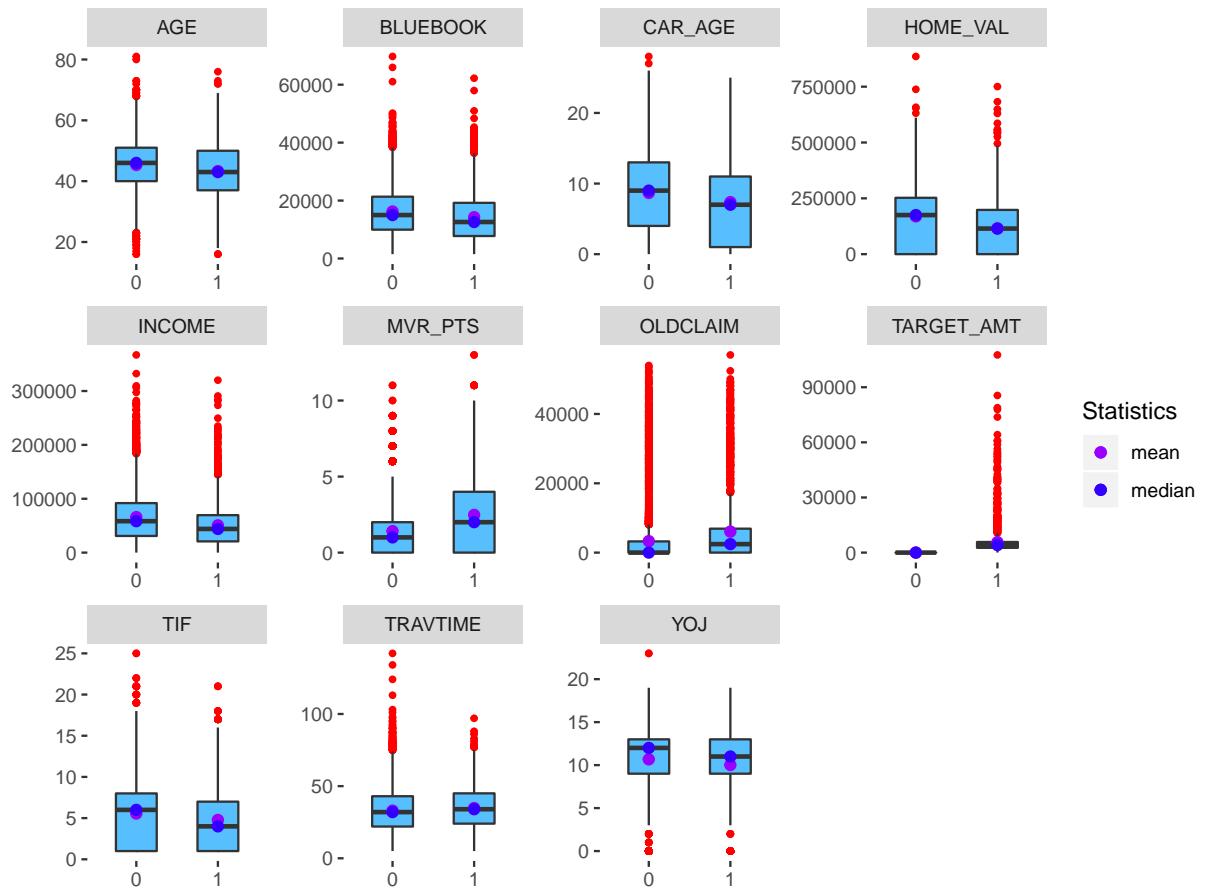


Figure 4: Linear relationship between each numeric predictor and the target

## 1.2 Linearity

[JO: DON'T THINK THE LOG TRANSFORM DISAMBIGUATES LINEAR RELATIONSHIPS - IF TEAM ALIGNED, WE CAN TRY TO ADD REGRESSION EQUATIONS TO THE FACETS: <https://community.rstudio.com/t/annotate-ggplot2-with-regression-equation-and-r-squared/6112/6>]

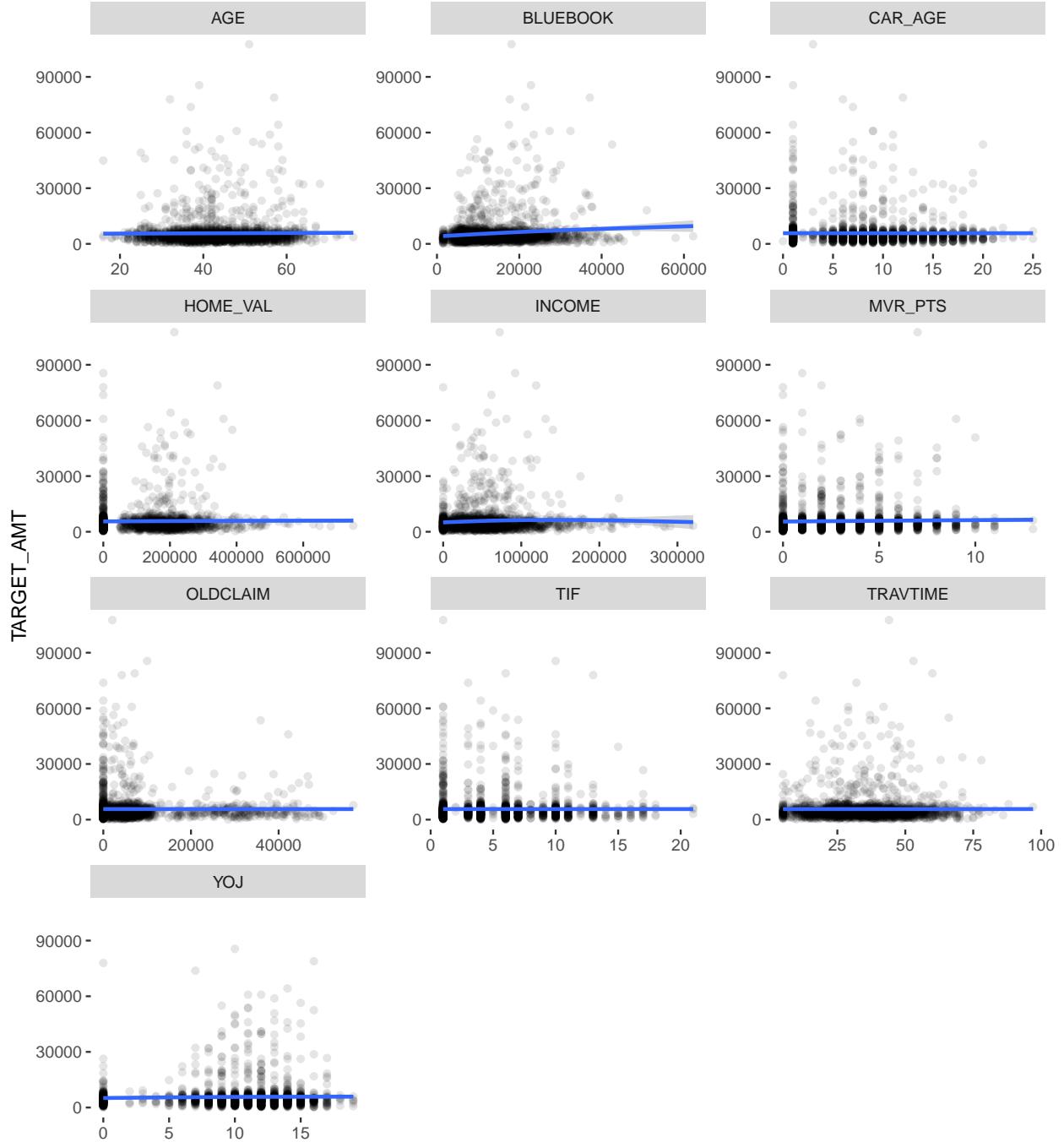


Figure 5: Scatter plot between numeric predictors and the TARGET\_AMT

The plotted numeric predictors with their raw values fail to show any clear linear relationships with the TARGET\_AMT except for the faintest of linearity in the BLUEBOOK variable.

[JO: DON'T THINK THE LOG TRANSFORM HAS HELPED - THINK WE SHOULD ADD]

[GAB: ADD WHAT? LOG DIDN'T HELP THOUGH YOU'RE RIGHT]

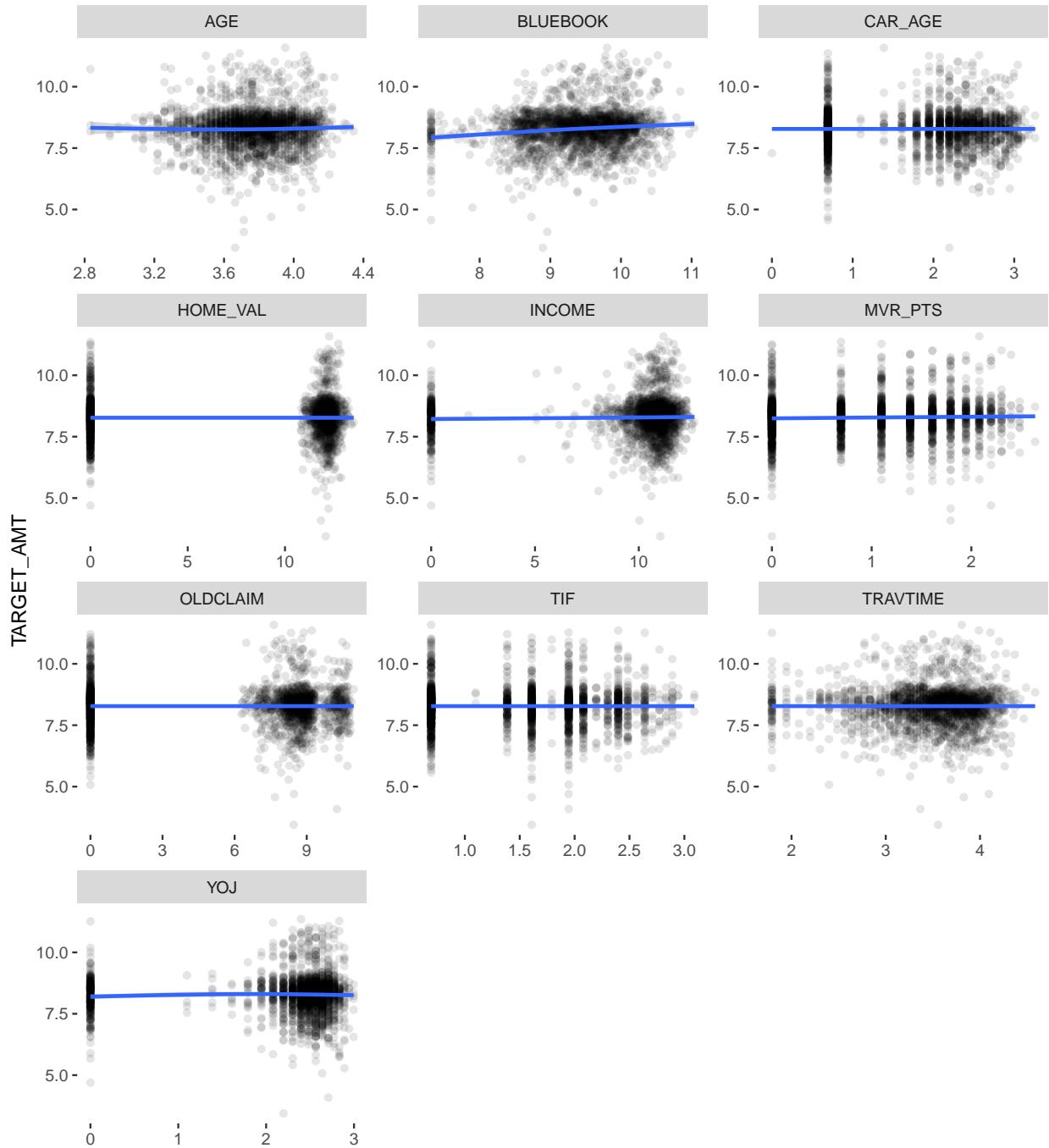


Figure 6: Scatter plot between log transformed numeric predictors and the log transformed TARGET\_AMT

In an attempt to distinguish the linearity of the variables alongside the TARGET\_AMT, all numeric predictors and the TARGET\_AMT underwent a log transformation. As a result, the linearity of BLUEBOOK became more apparent, but there was no obvious influence on the linearity of any of the other variables.

### 1.3 Missing Data

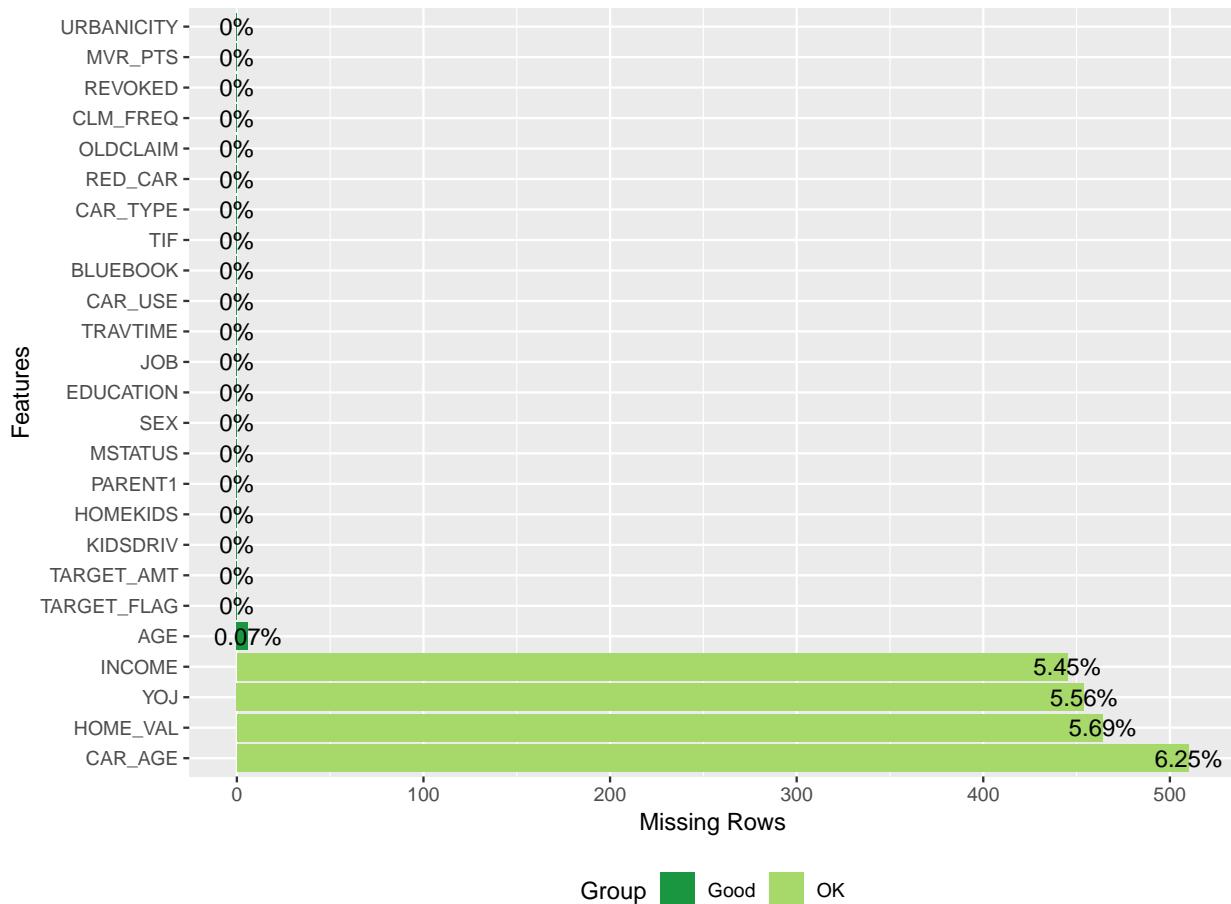


Figure 7: Missing data

A number of variables are missing observations: AGE, INCOME, YOJ, HOME\_VAL, CAR\_AGE. For AGE, the number is inconsequential, but the others range between 5% and 6% of total.

[JO: WHERE DID WE NET OUT ON CHECKING WHETHER THE SAME OBSERVATIONS ARE MISSING THE ABOVE VARIABLES AND WHETHER THEY HAVE ANYTHING IN COMMON I.E. THEY'RE ALL MARRIED, HIGH EARNERS, HOMEOWNERS?]

It was decided the missing values would be imputed when preparing the data for modeling.

## 2 DATA PREPARATION

### 2.1 Missing Values

[JO: THINK WE SHOULD DESCRIBE PURPOSE OF/ NEED FOR VALUE IMPUTATION. MICE IMPUTATION ASSUMES ‘MISSING AT RANDOM’ (MAR), SO THINK WE’LL NEED TO ESTABLISH THAT THIS IS THE CASE.][GAB: WE PROBABLY SHOULD BUT I THINK NO HARM NO FOUL IF WE SKIP IT]

[JO: WHAT’S THE DIFFERENCE BETWEEN THE M AND MAXIT VALUES (1 FOR AGE, 2 FOR OTHERS?)][GAB: CAN I KNOW THE ANSWER TO THIS TOO?]

To deal with missing data values for the variables `INCOME`, `YOJ`, `HOME_VAL`, and `CAR_AGE` - and to a lesser extent `AGE` - the MICE (Multivariate Imputation By Chained Equations) package was leveraged. The package assumes missing values are missing at random and creates multiple imputations (replacement values) for multivariate missing data using a method based on Fully Conditional Specification, where each incomplete variable is imputed by a separate model. The method can impute mixes of continuous, binary, unordered categorical and ordered categorical, and continuous two-level data; and it can maintain consistency between imputations by means of passive imputation. The quality of imputed values was inspected using multiple diagnostic plots.

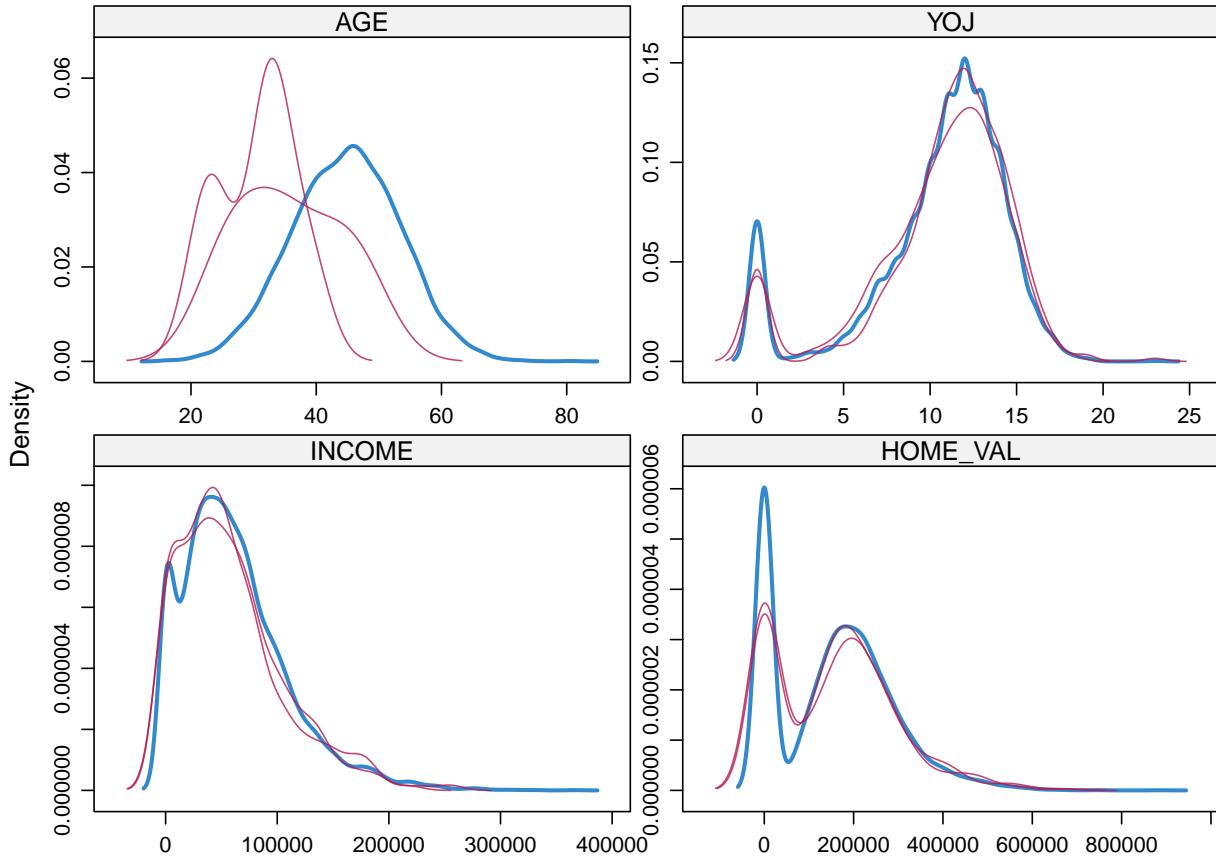
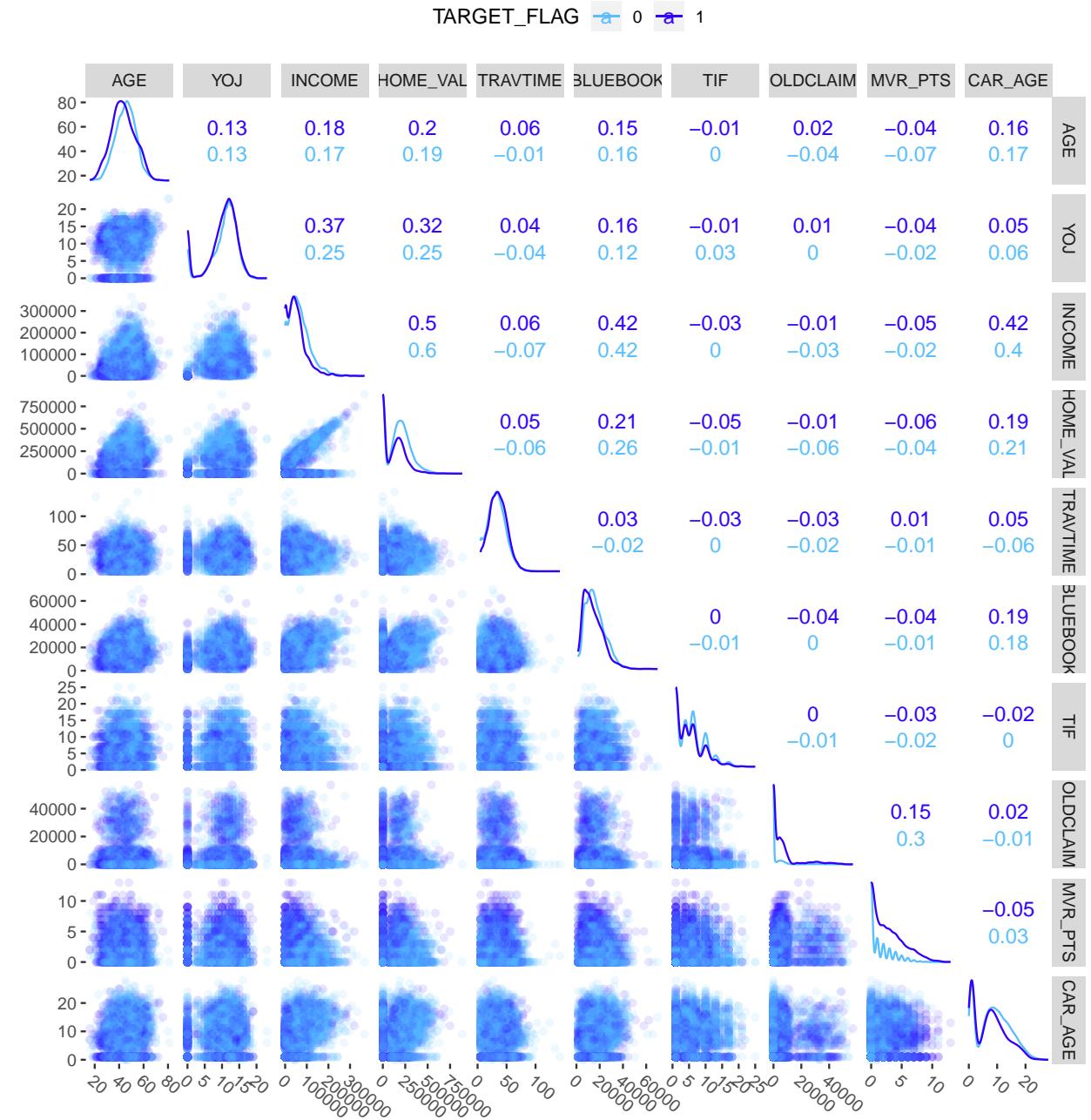


Figure 8: Difference between original and imputed data

The blue and red lines represent the distribution originally known and imputed values respectively. With the exception of AGE, the distribution of the imputed values accords with the distribution of pre-existing values. The imputed values will be used for the four variables. When it comes to AGE (only .07 or XX number of cases), it was to be imputed separately using median imputation.



Unsurprisingly, higher levels of INCOME are found with higher values of YOJ; this also means more income is disposable, which shows correlation with HOME\_VAL and BLUEBOOK.

Additionally, MVR\_PTS shows a positive correlation with OLDCLAIMS.



## 3 BUILD MODELS

### 3.1 Classification Models: Models 1, 2, 3, 4

The first four models take categorical variables as inputs and interpret their contributions to predicting the likelihood of a claim [JO: CONFIRM JUST ‘for new customers’?]. We use `drop` and `MASS::stepAIC` functions to judge which variables to remove, evaluating AIC statistics as we go.

#### 3.1.1 Model 1

```
TARGET_FLAG ~ NumParents + Male + EDUCATION + JOB + CAR_TYPE + RED_CAR  
+ REVOKED + Urban + Single + Commercial
```

For an easily interpretable model aimed at predicting TARGET\_FLAG, inputs for Model 1 were restricted to categorical variables alone. The AIC metric suggests that the `RED_CAR` variable can be removed.

#### 3.1.2 Model 2

```
TARGET_FLAG ~ KIDSDRV + AGE + HOMEKIDS + YOJ + INCOME + HOME_VAL  
+ TRAVTIME + BLUEBOOK + TIF + OLDCLAIM + CLM_FREQ + MVR PTS +  
CAR_AGE + PARENT1 + SEX + EDUCATION + JOB + CAR_TYPE + REVOKED +  
URBANICITY + MSTATUS + CAR_USE
```

As suggested by Model 1’s findings, Model 2 excludes the `RED_CAR` variable. AIC metrics suggest removing `AGE`, `CAR_AGE` and `SEX`.

#### 3.1.3 Model 3

```
TARGET_FLAG ~ KIDSDRV + HOMEKIDS + YOJ + INCOME + HOME_VAL + TRAV-  
TIME + BLUEBOOK + TIF + OLDCLAIM + CLM_FREQ + MVR PTS + PARENT1 + ED-  
UCATION + JOB + CAR_TYPE + REVOKED + URBANICITY + MSTATUS + CAR_USE
```

Building on model 2, model 3 excludes `AGE`, `CAR_AGE` and `SEX`.

[RJ: Anybody wants to create classification evaluation table in the select model section for those? – might want to roll back for the data cleaning I had done previously as we need the factorized values to create that matrices.]

#### 3.1.4 Model 4 - Binary logistic model

Model 4 incorporates all explanatory variables plus log transformations of skewed variables (`INCOME`, `TRAVTIME`, `BLUEBOOK`, `OLDCLAIM`, and `AGE`) in a binary logistic model refined through backward elimination.

Observations	7651
Dependent variable	TARGET_FLAG
Type	Generalized linear model
Family	binomial
Link	logit

$\chi^2(35)$	2010.61
Pseudo-R <sup>2</sup> (Cragg-Uhler)	0.34
Pseudo-R <sup>2</sup> (McFadden)	0.23
AIC	6875.46
BIC	7125.40

	Est.	S.E.	z val.	p	VIF
(Intercept)	1.99	0.82	2.43	0.02	NA
KIDSDRV1	0.55	0.11	5.01	0.00	1.13
KIDSDRV2	0.82	0.16	5.23	0.00	1.13
KIDSDRV3	0.83	0.32	2.64	0.01	1.13
KIDSDRV4	-11.59	196.87	-0.06	0.95	1.13
log(AGE)	-0.27	0.16	-1.69	0.09	1.26
YOJ	0.02	0.01	1.82	0.07	2.38
log(INCOME + 0.0000000000000001)	-0.02	0.00	-4.50	0.00	3.30
HOME_VAL	-0.00	0.00	-4.91	0.00	1.75
log(TRAVTIME)	0.41	0.05	7.67	0.00	1.03
log(BLUEBOOK)	-0.34	0.06	-5.91	0.00	1.49
TIF	-0.05	0.01	-7.11	0.00	1.01
log(OLDCLAIM + 0.0000000000000001)	0.01	0.00	6.50	0.00	1.26
MVR_PTS	0.10	0.01	6.86	0.00	1.25
PARENT1Yes	0.38	0.10	3.63	0.00	1.63
EDUCATIONBachelors	-0.41	0.11	-3.68	0.00	7.60
EDUCATIONMasters	-0.30	0.17	-1.80	0.07	7.60
EDUCATIONPhD	-0.28	0.20	-1.41	0.16	7.60
EDUCATIONHigh School	0.04	0.10	0.43	0.67	7.60
JOBClerical	0.50	0.20	2.48	0.01	27.64
JOBDoctor	-0.32	0.27	-1.18	0.24	27.64
JOBHome Maker	0.13	0.22	0.61	0.54	27.64
JOBLawyer	0.15	0.17	0.88	0.38	27.64
JOBManager	-0.58	0.18	-3.27	0.00	27.64
JOBProfessional	0.20	0.18	1.09	0.28	27.64
JOBStudent	0.05	0.23	0.22	0.83	27.64
JOBBlue Collar	0.38	0.19	2.00	0.04	27.64
CAR_TYPEPanel Truck	0.60	0.15	4.02	0.00	2.34
CAR_TYPEPickup	0.56	0.10	5.38	0.00	2.34
CAR_TYPESports Car	0.92	0.11	8.24	0.00	2.34
CAR_TYPEVan	0.63	0.13	5.00	0.00	2.34
CAR_TYPESUV	0.71	0.09	8.01	0.00	2.34
REVOKEDEYes	0.72	0.08	8.64	0.00	1.01
URBANICITYRural	-2.34	0.12	-20.07	0.00	1.14
MSTATUSNo	0.47	0.08	5.53	0.00	1.96
CAR_USEPrivate	-0.75	0.09	-7.91	0.00	2.45

Standard errors: MLE

### 3.2 Regression Model: Models 5, 6

The next two models are multiple linear regression models aimed at predicting the value of claims based on different approaches, including constraining the cases based on TARGET\_FLAG (i.e. based on whether or not

a claim was filed) and different approaches to selecting explanatory variables.

23 lines were removed where TARGET\_AMT greater than \$45,000; these lines had a BLUEBOOK value far less than the car crash cost. A new variable milage was created based on TRAVTIME and CAR\_AGE.

### 3.2.1 Model 5 - Multiple linear regression model

Model 5 is a multiple linear regression model built only on cases with claims where TARGET\_FLAG equals 1. The model is refined using stepwise elimination. From the model summary it can be observed that the Adjusted R-squared value is very low at 0.04.

Observations	1988
Dependent variable	TARGET_AMT
Type	OLS linear regression
<hr/>	
F(49,1938)	1.46
R <sup>2</sup>	0.04
Adj. R <sup>2</sup>	0.01

### 3.2.2 Model 6 - Multiple linear regression model

Model 6 is a multiple linear regression model built on all cases - in other words, it relaxed the constraint that a claim was filed, and so includes TARGET\_AMT values of 0. Forward elimination was used to refine variable selection. The Adjusted R-squared value significantly improved compared to the previous model.

## 4 SELECT MODELS

## 5 Appendix

The appendix is available as script.R file in project4\_insurance folder.

[https://github.com/betsyrosalen/DATA\\_621\\_Business\\_Analyt\\_and\\_Data\\_Mining](https://github.com/betsyrosalen/DATA_621_Business_Analyt_and_Data_Mining)

	Est.	S.E.	t val.	p
(Intercept)	13105.50	14676.12	0.89	0.37
KIDSDRV1	106.21	442.21	0.24	0.81
KIDSDRV2	425.38	583.12	0.73	0.47
KIDSDRV3	-257.74	1054.66	-0.24	0.81
log(AGE)	-4711.83	3445.76	-1.37	0.17
AGE	124.74	82.67	1.51	0.13
HOMEKIDS1	319.95	476.22	0.67	0.50
HOMEKIDS2	296.74	468.87	0.63	0.53
HOMEKIDS3	432.90	529.25	0.82	0.41
HOMEKIDS4	1207.11	821.24	1.47	0.14
HOMEKIDS5	942.71	2604.73	0.36	0.72
YOJ	-48.97	46.55	-1.05	0.29
log(INCOME + 0.0000000000000001)	19.03	17.37	1.10	0.27
INCOME	-0.01	0.00	-1.67	0.10
CAR_AGE	-68.70	35.23	-1.95	0.05
log(milage)	111.30	79.38	1.40	0.16
log(BLUEBOOK)	991.60	472.70	2.10	0.04
BLUEBOOK	-0.01	0.04	-0.20	0.84
TIF	-24.96	29.33	-0.85	0.39
log(OLDCLAIM + 0.0000000000000001)	188.24	331.78	0.57	0.57
OLDCLAIM	0.01	0.03	0.40	0.69
CLM_FREQ1	-7805.75	13351.39	-0.58	0.56
CLM_FREQ2	-7745.34	13348.65	-0.58	0.56
CLM_FREQ3	-8067.59	13354.06	-0.60	0.55
CLM_FREQ4	-9164.28	13353.88	-0.69	0.49
CLM_FREQ5	-8497.12	13546.44	-0.63	0.53
MVR_PTS	69.29	48.96	1.42	0.16
PARENT1Yes	-192.34	463.54	-0.41	0.68
SEXF	-429.18	416.33	-1.03	0.30
EDUCATIONBachelor	295.14	445.00	0.66	0.51
EDUCATIONMaster	1180.38	755.96	1.56	0.12
EDUCATIONPhD	2571.09	918.26	2.80	0.01
EDUCATIONHigh School	-62.62	353.00	-0.18	0.86
JOBClerical	163.08	829.16	0.20	0.84
JOBDoctor	-1356.27	1187.50	-1.14	0.25
JOBHome Maker	-38.33	897.73	-0.04	0.97
JOBLawyer	424.71	701.95	0.61	0.55
JOBManager	-308.89	736.63	-0.42	0.68
JOBProfessional	126.37	774.22	0.16	0.87
JOBStudent	26.74	907.40	0.03	0.98
JOBBlue Collar	775.16	788.77	0.98	0.33
CAR_TYPEPanel Truck	277.72	684.21	0.41	0.68
CAR_TYPEPickup	419.01	413.46	1.01	0.31
CAR_TYPESports Car	599.78	519.48	1.15	0.25
CAR_TYPEVan	-47.23	532.48	-0.09	0.93
CAR_TYPESUV	638.25	462.22	1.38	0.17
REVOKEDEYes	-603.17	361.54	-1.67	0.10
URBANICITYRural	-432.38	521.59	-0.83	0.41
MSTATUSNo	756.94	321.79	2.35	0.02
CAR_USEPrivate	-67.81	360.16	-0.19	0.85

Observations	7628
Dependent variable	TARGET_AMT
Type	OLS linear regression

F(51,7576)	112.78
R <sup>2</sup>	0.43
Adj. R <sup>2</sup>	0.43

	Est.	S.E.	t val.	p
(Intercept)	2191.67	4527.86	0.48	0.63
TARGET_FLAG	5062.37	77.46	65.35	0.00
KIDSDRV1	57.56	128.25	0.45	0.65
KIDSDRV2	216.37	181.01	1.20	0.23
KIDSDRV3	-78.46	359.75	-0.22	0.83
KIDSDRV4	-616.44	1857.51	-0.33	0.74
log(AGE)	-1396.62	1071.81	-1.30	0.19
AGE	36.01	24.97	1.44	0.15
HOMEKIDS1	45.75	123.87	0.37	0.71
HOMEKIDS2	-11.80	121.21	-0.10	0.92
HOMEKIDS3	42.75	142.17	0.30	0.76
HOMEKIDS4	329.08	235.37	1.40	0.16
HOMEKIDS5	353.06	758.81	0.47	0.64
YOJ	-10.11	11.22	-0.90	0.37
log(INCOME + 0.0000000000000001)	6.14	4.47	1.37	0.17
INCOME	-0.00	0.00	-1.66	0.10
CAR_AGE	-20.00	9.91	-2.02	0.04
log(milage)	48.09	32.61	1.47	0.14
log(BLUEBOOK)	359.64	126.73	2.84	0.00
BLUEBOOK	-0.01	0.01	-0.93	0.35
TIF	-4.70	7.24	-0.65	0.52
log(OLDCLAIM + 0.0000000000000001)	68.11	102.72	0.66	0.51
OLDCLAIM	-0.00	0.01	-0.04	0.97
CLM_FREQ1	-2816.92	4132.78	-0.68	0.50
CLM_FREQ2	-2790.20	4129.68	-0.68	0.50
CLM_FREQ3	-2913.99	4131.39	-0.71	0.48
CLM_FREQ4	-3347.18	4139.71	-0.81	0.42
CLM_FREQ5	-3187.07	4191.07	-0.76	0.45
MVR PTS	33.53	15.99	2.10	0.04
PARENT1Yes	34.95	128.82	0.27	0.79
SEXF	-93.04	96.57	-0.96	0.34
EDUCATIONBachelors	66.15	121.79	0.54	0.59
EDUCATIONMasters	235.53	180.20	1.31	0.19
EDUCATIONPhD	483.03	214.94	2.25	0.02
EDUCATIONHigh School	-11.92	101.47	-0.12	0.91
JOBClerical	-50.35	203.12	-0.25	0.80
JOBDoctor	-253.93	242.30	-1.05	0.29
JOBHome Maker	-26.14	223.12	-0.12	0.91
JOBLawyer	85.76	175.67	0.49	0.63
JOBManager	-61.15	171.65	-0.36	0.72
JOBProfessional	-38.28	183.24	-0.21	0.83
JOBStudent	-94.47	227.54	-0.42	0.68
JOBBlue Collar	142.24	191.48	0.74	0.46
CAR_TYPEPanel Truck	130.49	168.40	0.77	0.44
CAR_TYPEPickup	103.81	101.14	1.03	0.30
CAR_TYPESports Car	149.94	130.06	1.15	0.25
CAR_TYPEVan	13.69	126.08	0.11	0.91
CAR_TYPESUV	132.04	106.64	1.24	0.22
REVOKEDYes	-132.70	104.51	-1.27	0.20
URBANICITYRural	-22.84	85.38	-0.27	0.79
MSTATUSNo	21	161.54	2.10	0.04
CAR_USEPrivate		0.53	97.43	0.01

Standard errors: OLS